

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 06-202882

(43)Date of publication of application : 22.07.1994

(51)Int.Cl.

G06F 9/46

(21)Application number : 05-238548

(71)Applicant : INTERNATL BUSINESS MACH CORP
<IBM>

(22)Date of filing : 24.09.1993

(72)Inventor : CHEN MING-SYAN
TUREK JOHN J E
WOLF JOEL L
YU PHILIP SHI-LUNG

(30)Priority

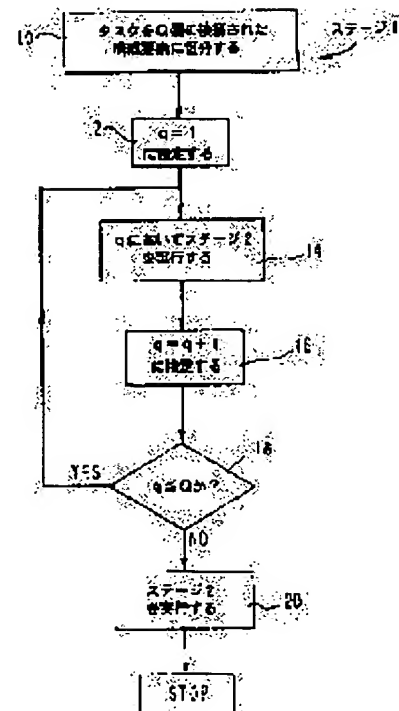
Priority number : 92 968717 Priority date : 30.10.1992 Priority country : US

(54) SCHEDULING METHOD FOR A PLURALITY OF TASKS AND JOBS

(57)Abstract:

PURPOSE: To provide an efficient and effective solution method for the flexible scheduling problems of tasks having priority.

CONSTITUTION: This scheduling method for a plurality of tasks and jobs consists of stages 1 to 3. The stage 1 is preliminary and divides the tasks into connected components (10). A prescribed job is calculated (14) at the stage 2 and each schedule of tasks contained in the job is decided according to its priority constraint. The stage 2 is repeated to calculate the execution time of jobs. The schedule of jobs having no priority constraint are decided (20) at the stage 3.



(19)日本国特許庁(JP)

(12)公開特許公報(A)

(11)特許出願公開番号

特開平6-202882

(43)公開日 平成6年(1994)7月22日

(51)Int.Cl.⁵

G 0 6 F 9/46

識別記号

3 4 0 B 8120-5B

庁内整理番号

F I

技術表示箇所

審査請求 有 請求項の数 7 (全 13 頁)

(21)出願番号 特願平5-238548

(22)出願日 平成5年(1993)9月24日

(31)優先権主張番号 9 6 8 7 1 7

(32)優先日 1992年10月30日

(33)優先権主張国 米国(US)

(71)出願人 390009531

インターナショナル・ビジネス・マシー
ズ・コーポレーション

INTERNATIONAL BUSIN
ESS MACHINES CORPO
RATION

アメリカ合衆国10504、ニューヨーク州
アーモンク (番地なし)

(72)発明者 ミンーシャン チェン

アメリカ合衆国10598、ニューヨーク州ヨ
ークタウン ハイツ、ブレンダー レイン
710

(74)代理人 弁理士 合田 深 (外6名)

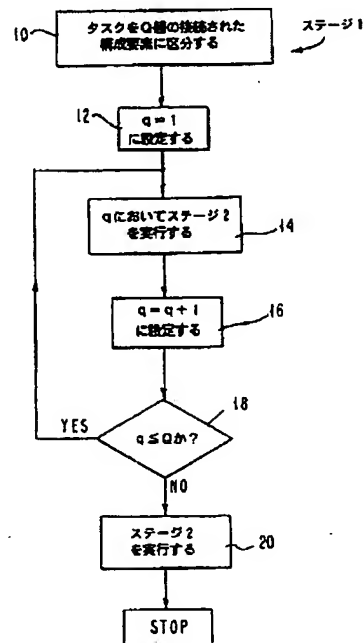
最終頁に続く

(54)【発明の名称】 複数のタスク及びジョブのスケジューリング方法

(57)【要約】

【目的】 優先順位を備えたタスクにおいて柔軟なスケジューリングの問題の効率的かつ効果的な解決法を提供する。

【構成】 複数のタスク及びジョブのスケジューリング方法は、ステージ1、2及び3から成る。ステージ1は予備ステージで、接続された構成要素にタスクを区分する(10)。ステージ2で、所定のジョブの計算を行い(14)、ジョブ内のタスクの各々のスケジュールを優先順位制約に従って決定する。このステージを繰り返してジョブ実行時間の計算を行う。ステージ3では、優先順位制約のないジョブのスケジュールを決定する(20)。



【特許請求の範囲】

【請求項1】 並列に作動する複数のプロセッサに優先順位制約を有する複数のタスクをスケジューリングするための方法であって、

(a) 複数のジョブを定義するステップを備え、前記ジョブの各々が前記タスクの一部と単一のジョブ内に含まれる前記タスクの対のみに関連する全ての優先順位の制約を含み、

(b) 各前記ジョブにおいて、前記ジョブの前記タスクに複数のタスクスケジューリングを作成するステップを備え、前記タスクスケジューリングの各々が前記ジョブに割り当てられる可能性のある異なる数の前記プロセッサに対応し、前記ジョブの前記タスク間でのあらゆる優先順位制約を尊重し、

(c) 前記タスクスケジューリングの各々に推定されたジョブ実行時間を決定するステップと、

(d) 前記ジョブの各々及び前記ジョブの各々に割り当てられる各異なる数のプロセッサに前記推定されたジョブ実行時間を用いて、前記ジョブの各々へのプロセッサの割り当てを決定するステップと、

(e) 前記決定された割り当てを用いて前記ジョブにジョブスケジューリングを作成するステップと、

(f) ステップ(e)で作成されたジョブスケジューリングを用いて前記プロセッサに前記ジョブを実行するステップと、

から成る複数のタスクのスケジューリング方法。

【請求項2】 ステップ(a)で定義された前記ジョブの各々がデータベース照会として定義される、請求項1に記載の複数のタスクのスケジューリング方法。

【請求項3】 ステップ(b)が前記ジョブの前記タスクの優先順位制約を尊重するように前記ジョブ内の前記タスクを順序付けるステップを含む、請求項1に記載の複数のタスクのスケジューリング方法。

【請求項4】 ステップ(d)が、前記タスクスケジューリングの各々毎に推定されたジョブ実行時間とそれに対応するプロセッサの数の積を計算するステップを含み、前記ジョブの各々へのプロセッサの最初の割り当てが前記ジョブの各々にそのように計算された積の内の最小の積に対応する、請求項1に記載の複数のタスクのスケジューリング方法。

【請求項5】 ステップ(e)で作成された前記ジョブスケジューリングが2次元ビンパッキング方法を用いて作成される、請求項1に記載の複数のタスクのスケジューリング方法。

【請求項6】 前記ジョブスケジューリングの前記ジョブのいずれかの内に遊休プロセッサの時間があるかどうかについてステップ(e)で作成された前記ジョブスケジューリングを解析し、そのような遊休プロセッサの時間があるならば、ステップ(d)において決定された前記割り当てを改良し、前記改良された割り当てに基づいて新たな

ジョブスケジューリングを作成するステップと、を更に含む請求項1に記載の複数のタスクのスケジューリング方法。

【請求項7】 並列に作動する複数のプロセッサを有する計算システムで複数の別個のジョブをスケジューリングするための方法であって、前記ジョブの各々が少なくとも1つのタスクから成り、少なくともいくつかの前記ジョブが優先順位制約のありうる複数のタスクから成り、前記プロセッサの各々が別個のタスクを同時に実行することが可能であり、

(a) 各ジョブにおいて推定されたジョブ実行時間のセットを作成するステップを備え、あらゆる前記セットの各前記推定されたジョブ実行時間が前記各ジョブの実行に専用可能な異なる数の前記プロセッサに対応し、前記推定されたジョブ実行時間が生成されると共に、前記各ジョブから成るタスクのあらゆる優先順位制約に従い、

(b) 各前記推定されたジョブ実行時間において、ジョブ実行時間と対応するプロセッサの数の積を計算するステップと、

(c) 推定されたジョブ実行時間の各前記セットにおいて、前記計算された積の内の最小の積を識別し、前記各ジョブの実行のために前記識別された最小の積に対応する前記プロセッサの数を試験的に割り当てるステップと、

(d) 各前記ジョブに前記試験的に割り当てられた数のプロセッサを用い、前記ジョブの各々の内の前記タスクに2次元ビンパッキングアルゴリズムを使用することによって、並列にある前記ジョブの全てを実行するための試験的な全体スケジューリングを作成するステップと、

(e) 前記試験的な全体スケジューリングに対応する全体実行時間を推定し、前記推定された全体実行時間を記録するステップと、

(f) 前記試験的な全体スケジューリングに遊休プロセッサがあるかどうかを決定するステップと、

(g) 遊休プロセッサがあるならば、識別されたジョブを実行するために試験的に割り当てられたプロセッサの最大数よりも少なく、遊休プロセッサの時間の最大量と対応する前記全体スケジューリングにおいてジョブを識別するステップと、

(h) 前記識別されたジョブを実行するために試験的に割り当てられたプロセッサの数を増加するステップと、

(i) 前記試験的な全体スケジューリングの前記遊休プロセッサの時間が、増加によってもはや減少できなくなるまでステップ(d)乃至(f)を繰り返すステップと、

(j) 前記記録された全体の実行時間から最小の全体実行時間を選択し、前記プロセッサで並列にある前記ジョブを実行するために前記最小の記録された全体実行時間に対応する前記試験的な全体スケジューリングを用いるステップと、

から成る複数のジョブのスケジューリング方法。

【発明の詳細な説明】

【0001】

【産業上の利用分野】本発明は、概して最後のタスクが終了する終了時間を最小限にするように並列処理コンピュータシステムで実行するための複数のタスクをスケジュールすることに関し、特に、タスク間に優先順位制約があるようなタスクをスケジュールすると共に、同時にスケジュール内の当該優先順位制約を尊重することに関する。例えば、タスクは複数のデータベース照会内で個々のステップであってもよい。

【0002】

【従来の技術】データベースからの情報要求（一般的に「データベース照会(database queries)」と称される）は概して数多くのステップ（本明細書では「タスク(task)」と称される）を含み、照会における少なくとも1つのタスクの結果（本明細書では「ジョブ(job)」と称される）は同一のジョブにおける後のタスクで用いられる。これは、ジョブ（もう1つのタスクで使用される結果を生じる）内のいかなる先のタスク（単数又は複数）も、当該結果を用いるいかなる後のタスクの実行が開始する前に終了されなければならない（即ち、「優先(precede)」しなければならない）ために、優先順位制約と称される。

【0003】並列処理システムは複数の処理要素を提供し、処理要素の各々は、他の処理要素が同一タスク若しくはジョブ又は異なるタスク若しくはジョブの他のタスク又は他のタスク部分を処理すると同時に、ジョブの全体のタスク又はそのようなタスクの一部を実行することができる。複数のジョブが並列処理システムで同時に実行されるとスケジュールされるとき、最小限の時間で全てのジョブにおいてタスクの総数の最後のタスクを終了するスケジュールを作成すると共に、勿論、タスク間の全ての優先順位制約を尊重することが一般的な目的となる。

【0004】本発明の目的のために、並列処理システムにおける各処理要素（以下本明細書で処理要素というよりもむしろ「プロセッサ(processor)」と称される）は、（単独又は少なくとも1つの他のプロセッサとの組合せで）一度に単一のタスクしか処理せず、もう1つのタスクを開始するために使用可能になる前に（単独又は少なくとも1つの他のプロセッサとの組合せで）タスクを処理して終了する。並列処理システムの物理資源が同時に1より多くのタスクを実行できるならば、プロセッサが1より多いと見られるだけである。

【0005】各ジョブはまた、割り当てられたプロセッサの数が非増加関数である実行時間を有すると仮定される。ジョブの実行時間は、ジョブの実行に割り当てられたプロセッサの数の増加と共に増加できない。これは、追加のプロセッサの使用(USE)によってジョブの終了時間が長くなる場合、最良のスケジュールはその追加され

たプロセッサを遊休(IDLE)にすることであり、従って、少なくともジョブの実行時間はもう1つのプロセッサの追加によって増加しないからである。

【0006】問題は、各ジョブへのプロセッサの割り当て（即ち、使用するプロセッサの数）、及び、全体として、最後のタスクが終了する終了時間を最小限にするように、全てのジョブのタスクをいかなる優先順位制約も尊重するプロセッサへ割り当てるスケジュールを見つけることである。

10 【0007】従って、P個のプロセッサから成るマルチプロセッサコンピュータシステム、並びに、このシステムでスケジュールされるN個のタスクのセットを考える。タスクは、優先順位関係

【0008】

【外1】

<

【0009】によって部分的に順序付けられる。各タスク $j \in \{1, \dots, N\}$ が任意の数のプロセッサ $\beta_i \in \{1, \dots, P\}$ に割り当てられ、タスク実行時間 $t_i(\beta_i) > 0$ が割り当てられたプロセッサの数の非増加関数であると仮定する。タスクに割り当てられた全てのプロセッサは、そのタスクを一致して実行することが必要とされる。すなわち、これらのプロセッサ β_i は全て任意の開始時間、例えば、 τ_i でタスク j を開始することが必要とされる。次にプロセッサは任意の後の終了時間 $\tau_i + t_i(\beta_i)$ でタスク j を終了する。各タスク $j \in \{1, \dots, N\}$ において、スケジュールはプロセッサ割り当て β_i 及び開始時間 τ_i から構成される。スケジュールは、以下の2つの意味で正当であることが必要とされる。・任意の時間 τ において、アクティブプロセッサの数はプロセッサの総数を越えない。換言すれば次の通りである。

【0010】

【数1】

$$\sum_{\{i | \tau_i \leq \tau < \tau_i + t_i(\beta_i)\}} \beta_i \leq P$$

【0011】・1つのタスク j_1 がもう1つのタスク j_2 に優先するならば、第2タスクは第1タスクが終了するまで開始できない。換言すれば次の通りである。

【0012】

【外2】

$$j_1 < j_2$$

【0013】ならば、 $\tau_{j_1} + t_{j_1}(\beta_{j_1}) \leq \tau_{j_2}$ である。最適なスケジュールを見つけることが問題であり、

【0014】

【数2】

$$\max_{1 \leq j \leq N} \{\tau_j + t_j(\beta_j)\}$$

【0015】によって与えられる全体のメイクスパン(makespan)が最小限にされる。換言すれば、最後のタスク終了時間を最小限にすることが目的である。これは柔軟な(malleable) スケジューリングの問題と称される。

【0016】図1は、3つのジョブに対する柔軟なスケジューリング問題を示すブロック図である。入力グラフのノード(節)によって表される3つのジョブから成るタスクと、エッジ(即ち、相互接続ライン)によって表される優先順位関係とから構成される。2つのノードがエッジによって接続されるならば、下位ノードに対応するタスクは上位ノードに対応するタスクが開始する前に終了しなければならない。スケジューリングプロシージャ(手順)は、出力として優先順位関係を尊重するタスクのスケジュールを作成する。プロセッサは横軸に、時間は縦軸に示されている。スケジュールのメイクスパンもまた示されている。

【0017】ドュー(J. Du)及びロイニング(J. Leung)著、「並列タスクシステムのスケジューリングの複雑さ(Complexity of Scheduling Parallel Task Systems)」(SIAM Journal of Discrete Mathematics (1989)は、柔軟なスケジューリングの問題は、優先順位制約がない特別の場合でさえも強い意味においてNP-hardであることを示している。(数学的に、これはその問題への全体的に最適な解決法を見つけるための効率的なプロシージャが発見される可能性は極めて低いことを意味する。)解決を試みるいかなるプロシージャも、この文献には展開されていない。

【0018】チェン(M.-S. Chen)、ユー(P. Yu)及びウー(K.-L. Wu)著、「マルチプロセッサシステムにおける複数の連結照会を実行するためのスケジューリング及びプロセッサの割り当て(Scheduling and Processor Allocation for the Execution of Multi-Join Queries in a Multiprocessor System)」(Proc. IEEE Conference on Data Engineering (1992)では、優先順位及び一般的なタスク実行時間を備えた単一の照会のコンテキストにおける柔軟なスケジューリングの問題が考慮されている。この文献では複数の照会の場合が考慮されておらず、単一の照会のために開発されたプロシージャは必ずしもあまり効果的ではない解決法を見つけている。

【0019】柔軟なスケジューリングの問題の非優先順位バージョンが、先の3つの文献で研究されている。これら文献の各々は、一般的なタスク実行時間を考慮し、少なくとも1つのプロシージャを提案していた。

【0020】チューレック(J. Turek)、ウルフ(J. Wolf)及びユー(P. Yu)著、「並列可能なタスクをスケジューリングするための近似アルゴリズム(Approximate Algor

ithms for Scheduling Parallelizable Tasks)」(Proc. Symposium on Parallel Algorithms and Architectures (1992)では、柔軟なプロシージャのクラスが提示され、プロシージャの各々が優先順位のない非柔軟なスケジューリングのためのプロシージャに基づき、プロシージャの各々が対応するより一層単純なプロシージャの最悪の場合の漸近の実行と一致する。しかしながら、優先順位はこの文献では考慮されていない。

【0021】チューレック(J. Turek)、ウルフ(J. Wolf)、パティパティ(K. Pattipati)及びユー(P. Yu)著、「並列可能なタスクのスケジューリング: 全てをシェルフに置く(Scheduling Parallelizable Tasks: Putting it all on the Shelf)」(Proc. ACM Sigmetrics Conference (1992)では、全ての可能なシェルフに基づく解決法のセットに対して最適なシェルフプロシージャが開発されている。(シェルフ解決法は、スケジューリング問題を解決するための方法のクラスを表す。)優先順位はこの文献でもまた考慮されていない。

【0022】クリシュナムルティ(R. Krishnamurthi)及びマー(E. Ma)著、「区分可能なマルチプロセッサシステムにおいて変化する区分サイズのタスクをスケジューリングするための近似アルゴリズム(An Approximation Algorithm for Scheduling Task on Varying Partition Sizes in Partitionable Multiprocessor Systems)」(IBM Research Report 15900 (1990)では、単一のシェルフに並列可能なタスクをバックする特別な場合を解決している。従って、この文献にも優先順位は考慮されていない。

【0023】

【発明が解決しようとする課題】従って、本発明の目的は、優先順位を備えたタスクにおいて柔軟なスケジューリングの問題の効率的かつ効果的な解決法を提供することである。

【0024】本発明のもう1つの目的は、各ジョブが優先順位制約の可能性のあるタスクの内の少なくとも1つのタスクから成る、複数のジョブを並列処理システムにスケジューリングするための方法を提供することである。

【0025】本発明の目的はまた、並列処理システムの複数のデータベース照会を効率的にスケジューリングするための方法を提供することである。

【0026】

【課題を解決するための手段と作用】これらの目的及び利点並びに追加の目的及び利点は、本発明によって達成される。本発明は、次の2つの内の第1番目のステージのみが優先順位制約を含む、スケジューリングの2つの別個のステージ(本明細書では、プロシージャの第2ステージ及び第3ステージと称される)を用いることによって、柔軟なスケジューリングの問題に対する効果的な解決法を決定するための効率的なプロシージャを提供する。

【0027】プロシージャの第1ステージ（予備でありしばしばオプションである）では、スケジュールされるべきタスクの全体のセットはジョブに区分され、各優先順位制約は単一のジョブにしか影響を及ぼさない。（即ち、特定のタスクが人力のために依存する全てのタスクは、その特定のタスクとして同一のジョブ内にも含まれる。）多くの場合、この条件を満たす適切なジョブは既に存在し、この第1ステージにおけるあらゆる実際の計算を不要にすることもある。これは、例えば、各照会が通常既に本質的にこの条件を満たしているデータベース照会の場合に一般的である。そのような場合には、各照会を別個のジョブとして定義するのが必要なだけである。

【0028】プロシージャの第2ステージ（実際のスケジューリングの第1ステージ）では、各ジョブは何度も別個で最適にスケジュールされ（そのジョブの実行に割り当てられる可能性のある異なる数の各プロセッサに一度）、各ジョブへのこれら他の最適なスケジュールの各々（本明細書では時々タスクスケジュールと称される）が作成されると共に、ジョブのタスク間での優先順位制約を尊重する。各ジョブにおけるこれら他の最適なスケジュールの各々の推定されたジョブ実行時間がまた決定され、このプロシージャの第3ステージで使用するために記録される。優先順位を尊重する最適なスケジュールを作成するためのあらゆる公知の方法が、この第2ステージで用いられる。「ダイナミックプログラミング(dynamic programming)」を用いて、この作成を行うための好ましい方法が開示されている。

【0029】プロシージャの第3ステージ（実際のスケジューリングの第2ステージ）では、このプロシージャの第2ステージで作成される各ジョブへの他の推定されたジョブ実行時間のセットは、入力として用いられて全てのジョブに最適な全体スケジュールを作成する。優先順位制約の各々が第2ステージ（即ち、第1スケジューリング段階）で既に完全に尊重されているため、このプロシージャの第3ステージは第2ステージで作成された推定されたジョブ実行時間からこの最適な全体スケジュールを作成するためにいかなる公知の方法を用いてもよい。ジョブのみがこのステージでスケジュールされるために、プロシージャの第3ステージでは尊重すべき優先順位制約はない。かつ、このプロシージャの第1ステージにおいて、ジョブ同士の間には優先順位制約がないようにジョブが定義されるために、ジョブ同士の間には優先順位制約はない。これは、この第3ステージであらゆるスケジューリング方法の使用も可能にするために故意に行われる。

【0030】このプロシージャの第3ステージにおいて、プロシージャの第2ステージで作成された推定されたジョブ実行時間の各セットの特定のメンバ（各ジョブに1セット）は、その選択に対応するプロセッサの数と

共に選択され、ジョブの全体スケジュールはこれらの選択を用いて作成される。概して、ジョブに対するプロセッサを割り当てるための多くの異なる選択も同様にスケジュールされる。全てのジョブの最小限の終了時間を有する全体スケジュールは、最適な全体スケジュールとして選択される。この選択を行う好ましい方法が記載されている。

【0031】記載された第3ステージの好ましい方法において、推定されたジョブ実行時間の各セットから、セットの効率的なメンバが選択され（即ち、対応する推定されたジョブ実行時間に関連して使用されるプロセッサの数が、そのタスクの実行に割り当てられたプロセッサの効率的な使用を反映するような意味において効率的に見えるもの、並びに、推定されたジョブ実行時間と使用されたプロセッサの数の積が最も小さいものが好ましい）、このメンバに対応するプロセッサの数が、そのジョブに割り当てられたプロセッサの最初の数として試験的に選択される。多くのプロセッサを各ジョブに対して割り当てるためのこの最初の試験的な選択を用いることによって、全体スケジュールはあらゆる公知の2次元ビンパッキング（詰め）アルゴリズムを用いて作成される。プロセッサを割り当てるためのこの最初のセットへの全てのジョブの推定された全体の実行時間が決定及び記録され、全体スケジュールで浪費された（即ち、遊休の）プロセッサの時間があるかどうかについて解析される。浪費されたプロセッサの時間があるならば、最も浪費されたプロセッサの時間と対応するジョブが識別され、そのジョブに割り当てられたプロセッサの数が1だけ増加され、新たな全体スケジュールが作成される。以前に記録された最良の時間より良ければ、この新たな全体スケジュールの全体の実行時間もまた推定及び記録される。新たな全体スケジュールもまた、浪費されたプロセッサの時間があるかどうかについて解析される。このプロセスは、あらゆる浪費されたプロセッサの時間がなくなるまで続き、最良の（即ち、最小の）全体実行時間を作成した全体スケジュールが最適なスケジュールとして選択されて、ジョブを実行するのに使用される。

【0032】第2ステージの好ましい方法において、ジョブのタスクはまずトポロジカル（構成）順に分類されるが、これはジョブのタスク間の優先順位制約の全てを尊重する順序である。トポロジカル順序は枝が相互に連結したタスクの木として表され、葉はジョブのタスクとして、枝の相互連結はタスク間の優先順位制約として表される。

【0033】次に、葉から始めて木の幹に向かって戻ると、各タスク（及びそのタスクの全ての子タスク）は何度も最適にスケジュールされ（即ち、そのタスク及び、もしあれば、その子タスクを実行するために使用されてもよい各可能な数のプロセッサにつき一度）、そのタスクに可能な他の最適なスケジュールのセット（そのタス

クに割り当てられてもよい各可能な数のプロセッサにつき1つの他の最適なスケジュール)を作成する。タスクが何度も交互に最適にスケジュールされつつあり(各異なるプロセッサの割り当ての可能性につき一度)、タスクが少なくとも1つの子タスクを有するときは常に、子タスクもまた親タスクと共に最適にスケジュールされる。これは優先順位制約を尊重するために実行される。このスケジュールを効率的に実行するために、最適なスケジュールと、子タスクがプロシーダの初期において個々に何度も交互に最適にスケジュールされたときに作成された子タスクの推定された実行時間とのセットが使用される。

【0034】タスクが直列と同様に並列で実行可能な2つの子タスクを有するならば(2つの子タスクの間に優先順位制約がないため)、各プロセッサに割り当てられる数の最適なスケジュールは、1つの可能性として互いに並列で実行され、他の可能性として互いに直列で実行され、勿論、互いに並列でスケジュールされるときには各子タスクに対するプロセッサの各可能な割り当てを分割して、それら子タスクが他のスケジュールを作成することによって決定される。

【0035】各タスクに関して、最適なスケジュールはそのタスクに割り当てられてもよい各可能な数のプロセッサで作成され、(そのような各最適なスケジュールの推定された実行時間と共に)これら最適なスケジュールが記録される。先に述べたように、子タスクの最適なスケジュール及び推定された実行時間のセットはその子タスクの親タスク(単数又は複数)を最適にスケジュールし、ジョブにおける最後の親タスク(即ち、ルートタスク)は、第3ステージのジョブで使用されるセットとなる最適なスケジュール及び対応する推定された実行時間のセットを作成する。

【0036】並列に作動する複数のプロセッサに優先順位制約を有する複数のタスクをスケジュールするための方法は、(a)複数のジョブを定義するステップを備え、前記ジョブの各々が前記タスクの一部と単一のジョブ内に含まれる前記タスクの対のみに関連する全ての優先順位制約を含み、(b)各前記ジョブにおいて、前記ジョブの前記タスクに複数のタスクスケジュールを作成するステップを備え、前記タスクスケジュールの各々が前記ジョブに割り当てられる可能性のある異なる数の前記プロセッサに対応し、前記ジョブの前記タスク間でのあらゆる優先順位制約を尊重し、(c)前記タスクスケジュールの各々に推定されたジョブ実行時間を決定するステップと、(d)前記ジョブの各々及び前記ジョブの各々に割り当てられる異なる数のプロセッサに前記推定されたジョブ実行時間を用いて、前記ジョブの各々へのプロセッサの割り当てを決定するステップと、

(e)前記決定された割り当てを用いて前記ジョブにジョブスケジュールを作成するステップと、(f)ステッ

ブ(e)で作成されたジョブスケジュールを用いて前記プロセッサに前記ジョブを実行するステップと、から成る。

【0037】並列に作動する複数のプロセッサを有する計算システムで複数の別個のジョブをスケジュールするための方法は、前記ジョブの各々が少なくとも1つのタスクから成り、少なくとも幾つかの前記ジョブが優先順位制約のありうる複数のタスクから成り、前記プロセッサの各々が別個のタスクを同時に実行することが可能であり、(a)各ジョブにおいて推定されたジョブ実行時間のセットを作成するステップを備え、あらゆる前記セットの各前記推定されたジョブ実行時間が前記各ジョブの実行に専用可能な異なる数の前記プロセッサに対応し、前記推定されたジョブ実行時間が生成されると共に、前記各ジョブから成るタスクのあらゆる優先順位制約に従い、(b)各前記推定されたジョブ実行時間において、ジョブ実行時間と対応するプロセッサの数の積を計算するステップと、(c)推定されたジョブ実行時間の各前記セットにおいて、前記計算された積の内の最小の積を識別し、前記各ジョブの実行のために前記識別された最小の積に対応する前記プロセッサの数を試験的に割り当てるステップと、(d)各前記ジョブに前記試験的に割り当てられた数のプロセッサを用い、前記ジョブの各々の内の前記タスクに2次元ビンパッキングアルゴリズムを使用することによって、並列にある前記ジョブの全てを実行するための試験的な全体スケジュールを作成するステップと、(e)前記試験的な全体スケジュールに対応する全体実行時間を推定し、前記推定された全体実行時間を記録するステップと、(f)前記試験的な全体スケジュールに遊休プロセッサがあるかどうかを決定するステップと、(g)遊休プロセッサがあるならば、識別されたジョブを実行するために試験的に割り当てられたプロセッサの最大数よりも少なく、遊休プロセッサの時間の最大量と対応する前記全体スケジュールにおいてジョブを識別するステップと、(h)前記識別されたジョブを実行するために試験的に割り当てられたプロセッサの数を増加するステップと、(i)前記試験的な全体スケジュールの前記遊休プロセッサの時間が、増加によってもはや減少できなくなるまでステップ(d)乃至(f)を繰り返すステップと、(j)前記記録された全体の実行時間から最小の全体実行時間を選択し、前記プロセッサで並列にある前記ジョブを実行するために前記最小の記録された全体実行時間に対応する前記試験的な全体スケジュールを用いるステップと、から成る。

【0038】

【実施例】以下の記載において、簡潔性のために、タスク同士の間の優先順位制約が森を形成すると仮定する。(数学的に、これはエッジのシーケンスを介して所定のノードへ接続される全てのノードのセットが木を形成することを意味する。このセットは、所定のノードの接続

された構成要素と称される。そのような木の集合を森と称する。)この仮定は緩めることができる。図2を特に参照すると、本発明のステージから成るブロック図が示されている。

【0039】ブロック10で示されるステージ1は、本発明の予備ステージである。ステージ1は接続された構成要素にタスクを区分して、 i と j の双方が同じ区画にある場合のみ、

【0040】

【外3】

$i < j$ である。

【0041】換言すれば、各 q は森の木に正確に対応する。 Q 個のそのような接続された構成要素があると仮定する。区分のためのプロシーダは、コーレム(T. Corem)、リーサーソン(C. Leiserson)、リベスト(R. Rivest)によるアルゴリズム(McGrawHill (1992), 441-442頁)に見られる。タスクの Q 個のセットの各々が単一のジョブから構成されるとして考えるのは好都合である。

(勿論、ジョブの識別子が予め知られている場合は、そのようなプロシーダは予備段階として必要とされない。)代表となるデータベースの例において、ジョブは個々の照会及び照会内のステップに対するタスクに対応する。照会の識別子は予め知られている可能性が高い。

【0042】本発明の残りは事実上階層的である。2つの追加のステージ、即ち、ステージ2及びステージ3がある。

【0043】1と Q の間の所定のジョブ q におけるステージ2の計算がブロック14に示されている。この計算はブロック12で開始され、ブロック16及び18によって制御される。1とプロセッサの総数 P の間にあるプロセッサの各数 p において、ステージ2はジョブ q 内のタスクの各々のスケジュールを決定する。このスケジュールは関連する優先順位制約に従う。このステージを繰り返し用いることによって、各 q 及び各 p のジョブ実行時間の計算が行われる。ステージ2のより詳細な記述は以下に示される。

【0044】ステージ3はブロック20に示されている。このステージは、ジョブ同士の間には優先順位制約のない Q 個のジョブでスケジュールを決定する。ステージ3はステージ2の出力を入力として使用する。ステージ3のより詳細な記述は以下に示される。

【0045】ステージ2の詳細は図3のブロック図に示されている。ステージ2のプロシーダの性質を理解するために、図4に示されるジョブ優先順位ツリーを考える。(ノードはタスクを表す。2つのノードがエッジによって接続されるならば、上位ノードに対応するタスクが開始する前に下位ノードに対応するタスクが終了しなければならない。)ステージ2は以下の特性に従うという意味で連続する改良されたプロシーダである。即ち、タスクが所定のプロセッサに割り当てられるなら

ば、全ての子(及び結果として全ての子孫全体)はそれらプロセッサのサブセットに割り当てられる。従って、ステージ2のプロシーダは非常に強い方法で優先順位ツリーの構造を尊重する。

【0046】図5は、図4に示されるジョブの可能なステージ2の解決法を示している。陰影の付いた領域は浪費された作業を表す。横軸はプロセッサを表し、縦軸は時間を表す。なお、ルートタスク(ノード12)は全てのプロセッサに割り当てられる。ノード12の2つの

子、即ち、ノード10及び11、はそれ自体がサブツリーのルートである。ノード10に対応するサブツリー全体は、並列として特徴付けられるようにノード11と対応するサブツリー全体の左側にバックされる。左側のサブツリーを降りながら、ノード10の子と対応するツリー、即ち、ノード7及び8のツリーを考える。ノード7と対応するサブツリーは、バッチとして特徴付けられるようにノード8と対応するサブツリーより上にバックされる。実際、ステージ2のプロシーダは、所定のノードの全ての子に対応するサブツリーが並列方法又はバッチ方法でバックされる特性に従う。

【0047】ステージ2のプロシーダは、上記2つの特性を尊重する最適なスケジュールを見つける。図3を再び参照する。特に、本発明のステージ2はタスクをトポロジカルに順序付けることによって開始するため、

【0048】

【外4】

$j_1 < j_2$

【0049】は $j_1 < j_2$ を意味する。これはブロック30に示されている。トポロジカル分類のプロシーダは、コーレム(T. Corem)、リーサーソン(C. Leiserson)、リベスト(R. Rivest)によるアルゴリズム(McGrawHill (1992), 485-488頁)に見られる。次に、ダイナミックプログラミングを用いて、葉のノードを上へ、より小さな数のプロセッサからより大きな数のプロセッサの順に従って進む。 p 個のプロセッサを使用するタスク j 及びその全ての子孫において、ステージ2のプロシーダによって見つけられる最適なメイクスパンを、

【0050】

【外5】

$i(p)$

【0051】として定義する(ここで、 $1 \leq p \leq P$ とする)。

【0052】

【外6】

$i(p)$

【0053】を生成するため、ステージ2のプロシーダは p 個のプロセッサを使用してタスク j をバックし、いずれも良い(a)バッチ方法、又は(b)最も可能な

13

並列方法でバックされた、jの子のサブツリーの全ての
バックをタスクjより下に追加する。

【0054】ブロック36は、所定のタスクj及びプロ
セッサの数pにおける

【0055】

【外7】

$$\hat{i}_j(p)$$

【0056】の計算を示す。pループの開始がブロック
32に示され、pループはブロック42及び44によつて
制御される。jループの開始がブロック34に示さ
れ、jループはブロック38及び40によって制御され
る。

【0057】

【外8】

$$\hat{i}_j(p)$$

*

$$\hat{i}_j(p) = i_j(p) + \min \{ \hat{i}_{ch_1}(p) + \hat{i}_{ch_2}(p), \min_{1 \leq q \leq p} \max \{ \hat{i}_{ch_1}(q), \hat{i}_{ch_2}(p-q) \} \}$$

【0064】に設定する等となる。ジョブqに対してス
テージ2のプロシージャを使用することで、1とPの間
のプロセッサの各数pにおいてジョブ実行時間

【0065】

【数6】

$$T_q(p) = \hat{i}_1(p)$$

【0066】が生じる。

【0067】ステージ3のプロシージャは、サブルーチ
ンとしてプロシージャSUBを用いる。このサブルーチ
ンにおいて、各ジョブqは固定数P。のプロセッサで実
行されるとみなされて、固定された時間の量、例えばT
。で実行する。サブルーチンは、優先順位制約のない
ようなタスクのスケジュールを見つける。SUBは、時
間の関数として使用されるプロセッサの数のヒストグラ
ムを利用する。当該ヒストグラムは図6の陰影の部分と
して示される。横軸は使用されるプロセッサの数に対応
する。縦軸は時間を表す。タスクが開始及び終了するよ
うスケジュールされるため、ヒストグラムの形状も変化
する。SUBプロシージャの詳細は、図7のブロック図
に示されている。最初に、ブロック50において、タス
クはタスク時間の最も長いものから最も短いものに、順
に分類される。従って、 $T_{q_1} > T_{q_2}$ ならばリストの q_1
の前に q_2 が現れる。ブロック52は、順序付けられた
実行可能リストの作成を表し、全てのタスクは最初に
{1, ..., Q}である。時間tの順序付けられたリストも
また維持される。所謂現在の時間tは、常に時間リスト
の最小の時間である。(時間リストは、実際、現在の時※

$$\hat{t} - t < T_q \text{ となるように } t \text{ における全ての } \hat{t} \text{ で } H(\hat{t}) \text{ を } H(\hat{t}) + p_q$$

【0070】に設定することによって更新される。その

14

* 【0058】を計算するため、 CH_i がjの子の数を示
すとする。 $CH_i = 0$ ならば、

【0059】

【数3】

$$\hat{i}_j(p) = i_j(p)$$

【0060】に設定する。 $CH_i = 1$ ならば、

【0061】

【数4】

$$\hat{i}_j(p) = i_j(p) + \hat{i}_{ch_1}(p)$$

【0062】に設定する。 $CH_i = 2$ ならば、

【0063】

【数5】

20※間以後の全てのタスク終了時間から成る。) リストは最
初に時間 $t=0$ のみから構成されるため、その時間がプロ
シージャの開始時における現在の時間である。これら
の初期設定は、ブロック54及び56に示される。勿
論、ブロック58に示されるように、最初のヒストグラ
ムも同様に0である。現在の時間はSUBの実行中には
決して減少しない。更に、タスクはプロシージャが実行
するときに、常に現在の時間に等しい開始時間でプロセ
ッサに割り当てられるので、タスクの開始時間は以前に
割り当てられたあらゆるタスクの開始時間と少なくとも
同じ位の大きさになる。

【0068】ブロック60に示されるように、一般的な
プロシージャのステップは実行可能リストが空になるま
で繰り返される。一般的なステップは以下に示される。
図6を考えると、現在の時間tが示されていると仮定す
る。この時点で使用可能な(即ち、未使用の)プロセッ
サの数がヒストグラムから読み出される。SUBは、プロ
セッサの数と同じ数を必要とするタスクの実行可能リ
ストを検索する。ブロック62において、そうしたタス
クがあるかどうかテストされる。そうしたタスクがあ
るならば、リストの最初のタスクqが選ばれ、ブロック
64でtに等しい開始時間を割り当てて、終了時間が
 $t + T_q$ になる。タスクqはブロック66でリストから取
り除かれ、タイム $t + T_q$ がブロック68の時間リスト
に追加され、ヒストグラムがブロック70において

【0069】

【外9】

50 ようなタスクがないならば、現在の時間tはブロック7

2において時間リストから取り除かれ、時間リストの次の最小の時間がブロック74において新たな現在の時間となる。ヒストグラムは、現在の時間以後は常に単調で増加しない。従って、タスクの開始時間でプロセッサの可用性をチェックするのは不十分である。タスクが開始可能ならば、同じプロセッサで干渉なしに終了までを実行することができる。

【0071】次に、ステージ3自体を示す。図8のブロック図を参照する。各ジョブqにおいて、全体の作業 p 、 (p) が最小限にされるようにプロセッサの数 p_q から開始し、 $T_q = t$ 、 (p_q) に設定される。これはブロック82で示される。qループの開始はブロック80で示され、ループはブロック84及び86によって制御される。ブロック88において、浪費された作業Wが非ゼロに設定され、日付(date)を付けることになった最良のメイクスパンMが無限大に設定される。ブロック90に示されるように、浪費された作業Wが0に駆動されるまで繰り返される。一般的なステップは以下に示される。ブロック92に示されるように、非柔軟なスケジューリングの問題につきSUBを呼び出す。ブロック94において、日付を付けることになった最良のメイクスパンMがMの以前の値の最小値と、SUBプロセスの実行によって見つけれられるメイクスパンとに修正される。関数 $P-H(\tau)$ は、各時間 τ における未使用のプロセッサの数を示す。各タスクqを、qが実行していた時間の間に生じた浪費された作業の量と対応することができる。すなわち、

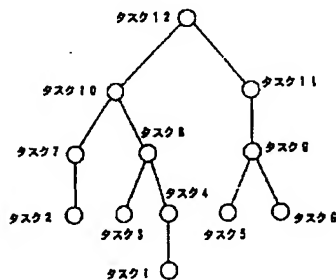
【0072】

【数7】

$$W_q = \int_{t_q}^{t_q+T_q} (P - H(\tau)) d\tau$$

【0073】に設定する。この計算はブロック98で示される。qループの開始はブロック96で示され、ループ

【図4】



* プはブロック100及び102で制御される。最大値W_qを備えたタスクをボトルネックタスクとみなし、ブロック104に示されるように浪費された作業Wをこの最大値として修正する。そのタスクqにおいて、プロセッサの数を、最小量の作業 $p T_q$ 、 (p) を備えた値 $p > p_q$ に増加する。次にブロック106に示されるように、 $p_q = p$ 、 $T_q = T_q$ 、 (p_q) に再設定する。ブロック92において再びSUBを読み出す。プロセスはボトルネックが0に駆動された固定数のステップで終了する。すなわち、全てのqで $W_q = 0$ になる。日付を付けるために得られる最良のメイクスパンが最終的な解決法である。

【0074】

【発明の効果】本発明は上記より構成され、優先順位を備えたタスクにおいて柔軟なスケジューリングの問題の効率的かつ効果的な解決法が提供される。

【図面の簡単な説明】

【図1】柔軟なスケジューリング問題の入力及び出力を示す高度なフローチャートである。

【図2】本発明に従ったプロセスの3つのステージの高度なフローチャートである。

【図3】図2のプロセスの第2ステージのフローチャートである。

【図4】ジョブ優先順位ツリーの一例である。

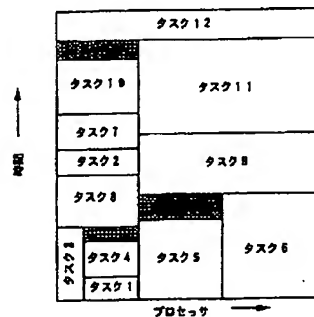
【図5】図4の例に対する第2ステージの解決法の一例である。

【図6】図2のプロセスの第3ステージにおいてSUBプロセスによって使用されるヒストグラムである。

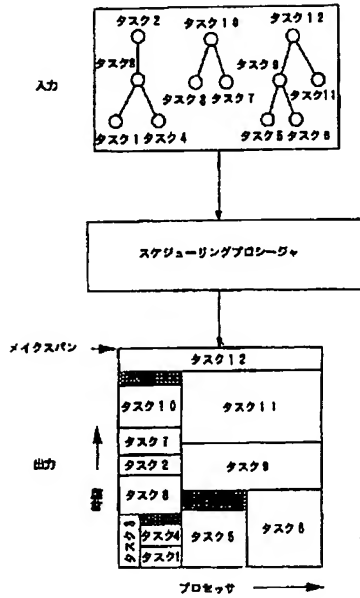
【図7】図2のプロセスの第3ステージにおけるSUBルーチンのフローチャートである。

【図8】図2のプロセスの第3ステージのフローチャートである。

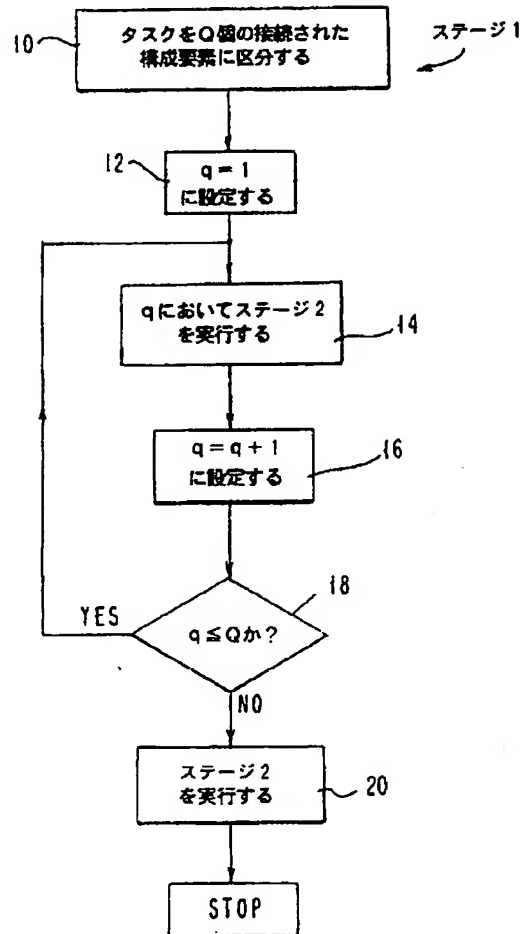
【図5】



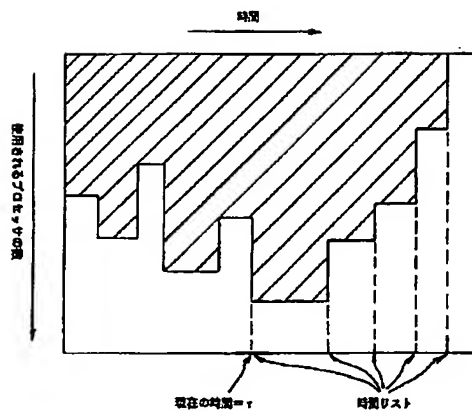
【図1】



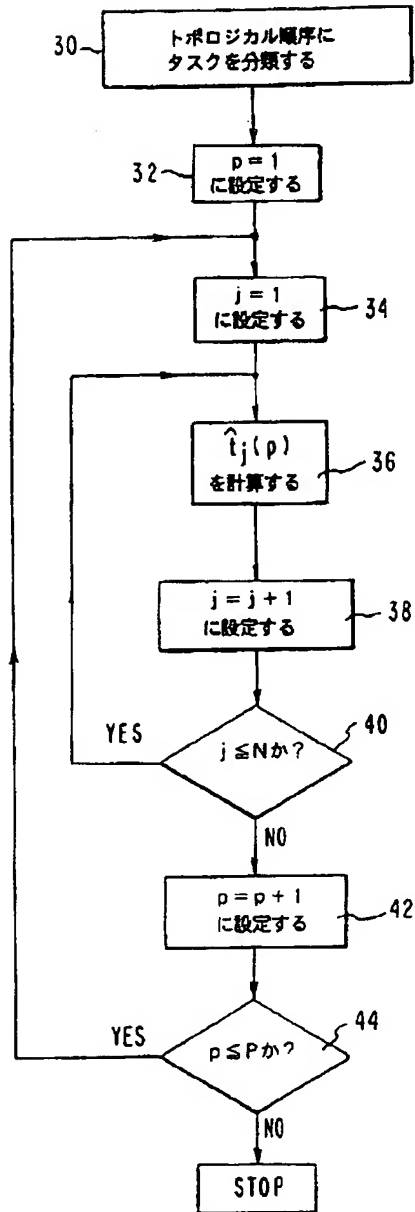
【図2】



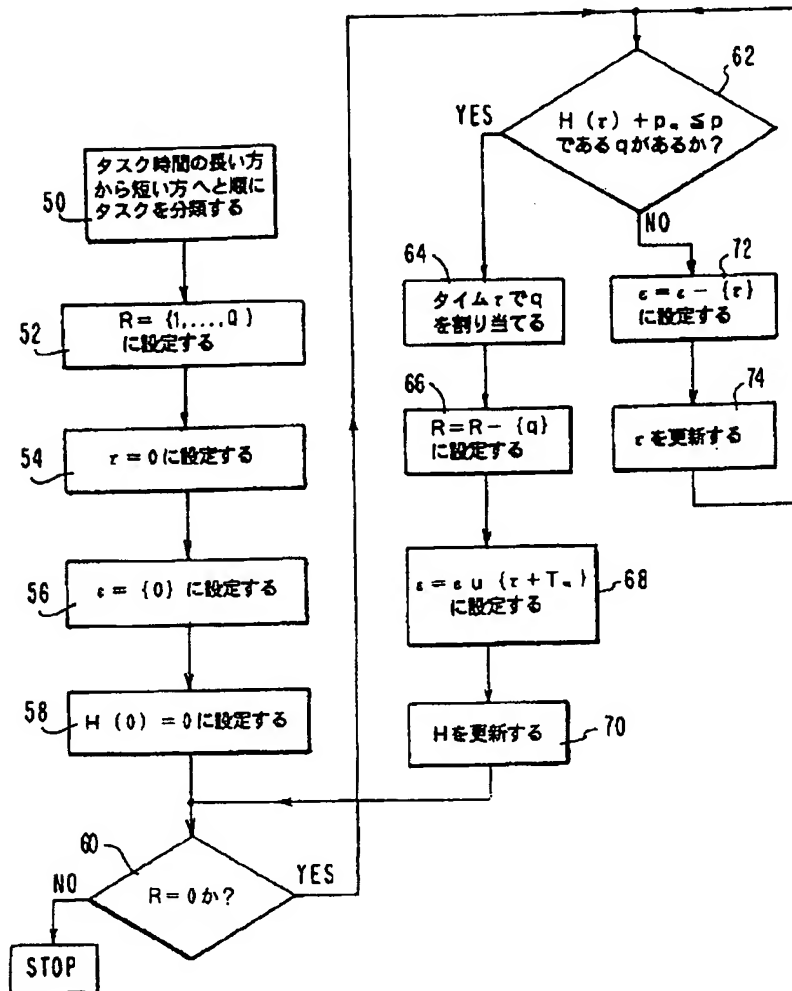
【図6】



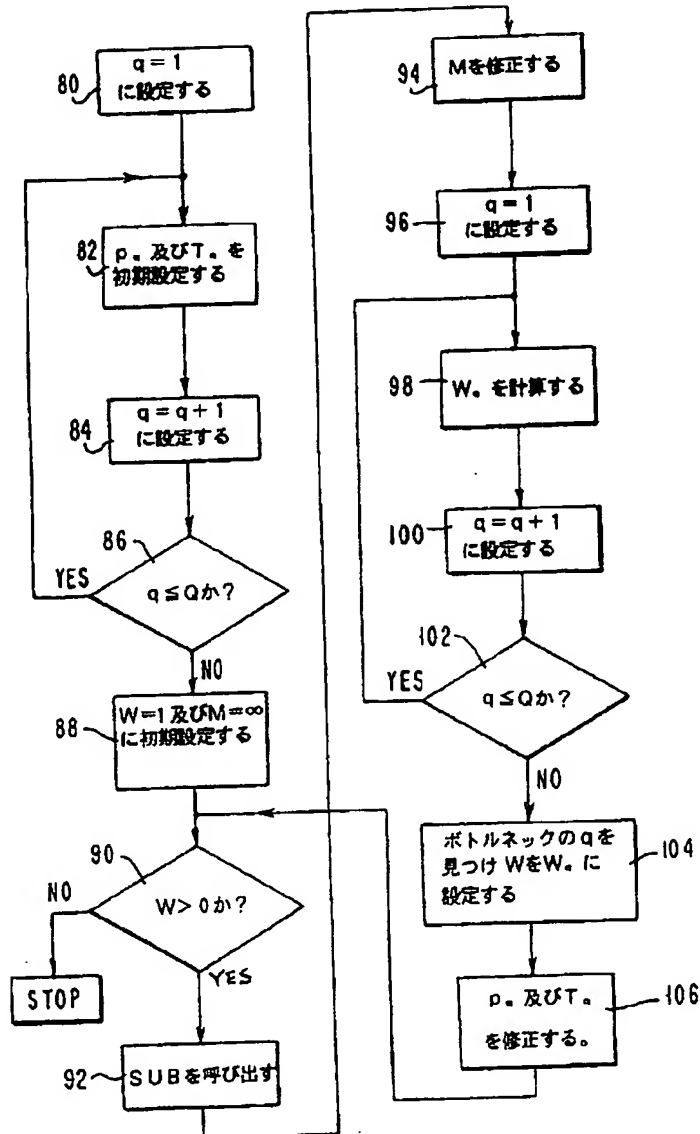
【図3】



【図7】



【図8】



フロントページの続き

(72)発明者 ジョン ジョゼフ エドワード トュレク
 アメリカ合衆国10514、ニューヨーク州チ
 ャッパクア、クロス リッジ ロード 50

(72)発明者 ジョエル レオナード ウォルフ
 アメリカ合衆国10536、ニューヨーク州カ
 トナー、チェロキー コート 7

(72)発明者 フィリップ シールン ユー
 アメリカ合衆国10514、ニューヨーク州チ
 ャッパクア、ストーンワイエ 18